

OPTIMIZING SPAM FILTERING WITH MACHINE LEARNING

TEAM ID: NM2023TMID17659

TEAM SIZE : 04

TEAM LEADER : GOVINDASAMY M (20UCS4712)

TEAM MEMBER : 1. GOVINDHARAJ B (20UCS4713)

2. VISHWANATHAN T (20UCS4741)

3. VIGNESH V (20UCS4740)

OPTIMIZING SPAM FILTERING WITH MACHINE LEARNING

1. INTRODUCTION :

Spam messages refer to unsolicited or unwanted messages/emails that are sent in bulk to users. In most messaging/emailing services, messages are detected as spam automatically so that these messages do not unnecessarily flood the users' inboxes. These messages are usually promotional and peculiar in nature. Thus, it is possible for us to build ML/DL models that can detect Spam messages.

OVERVIEW :

A spam filter is a program used to detect unsolicited, unwanted and virus-infected emails and prevent those messages from getting to a user's inbox. Like other types of filtering programs, a spam filter looks for specific criteria on which to base its judgments.

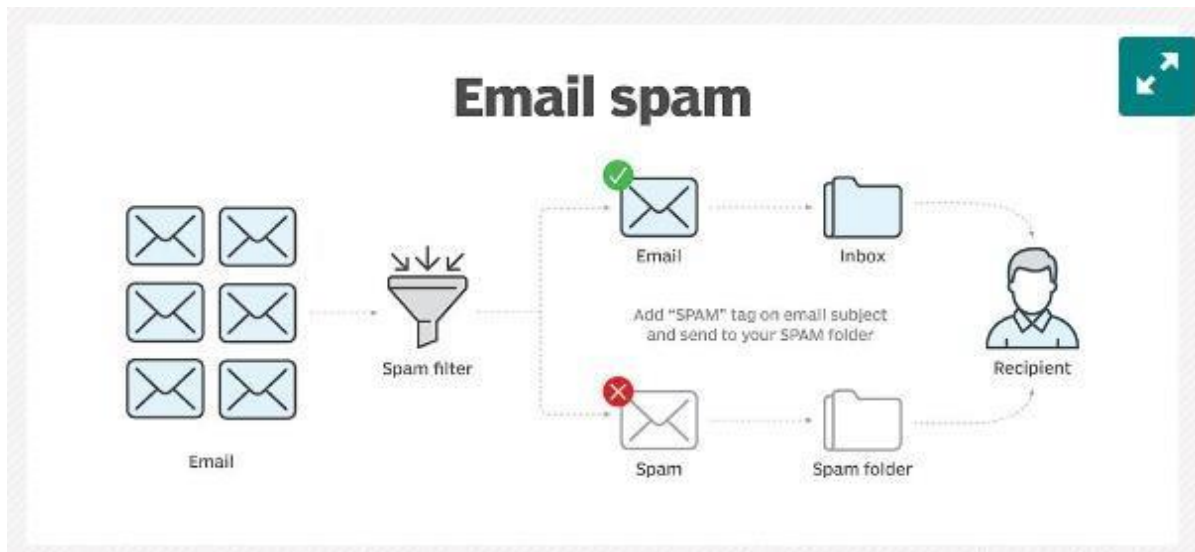
Internet service providers (ISPs), free online email services and businesses use email spam filtering tools to minimize the risk of distributing spam. For example, one of the simplest and earliest versions of spam filtering, like the one that was used by Microsoft's Hotmail, was set to watch out for particular words in the subject lines of messages. An email was excluded from the user's inbox whenever the filter recognized one of the specified words.

This method is not especially effective and often omits perfectly legitimate messages, called *false positives*, while letting actual spam messages through.

More sophisticated programs, such as Bayesian filters and other heuristic filters, identify spam messages by recognizing suspicious word patterns or word frequency. They do this by learning the user's preferences based on the emails marked as spam. The spam software then creates rules and applies them to future emails that target the user's inbox.

For example, whenever users mark emails from a specific sender as spam, the Bayesian filter recognizes the pattern and automatically moves future emails from that sender to the spam folder.

ISPs apply spam filters to both inbound and outbound emails. However, small to medium enterprises usually focus on inbound filters to protect their network. There are also many different spam filtering solutions available. They can be hosted in the cloud, hosted on servers or integrated into email software, such as Microsoft Outlook.



PURPOSE :

Spam filters are algorithms that detect unsolicited, undesired or infected emails, and block those messages from reaching inboxes. There are a variety of spam filters, with detection capabilities ranging from basic pattern matching, all the way through to machine learning.

Spam filters do a lot of work. And much of it is complex, confusing and constantly changing. With the pressure to keep viruses out and determine what's unsolicited vs solicited, it's no wonder they sometimes don't get everything right.

Don't worry though, there's lots you can do to ensure you're playing by the rules and feeding spam filters the information they need, to make an informed decision about your email. More on the steps you can take later.

There are three main types of spam filters: Email service provider, third-party and desktop.

1. Email service provider spam filters

These are the spam filters you're likely using without even knowing. Email service providers are the cloud based services that most people use to send, receive and organize their emails. Think Gmail, Yahoo Mail, Office 365 and Hotmail (yes, Hotmail is still a thing!).

While you may think of those names being mainly used by individuals, many of them offer services packaged for companies to use — en masse with their employees — with one popular service being Google Workspace.

These email service providers have their own built-in spam filters with advanced algorithms that look for the latest spam patterns and apply machine learning (ML) to evaluate a sender's trust and the message's inbox-worthiness.

2. Third-party spam filters

Third party spam filters are additional layers of spam protection that companies add to their email system (whether it's onsite or hosted in the cloud). They work hand-in-hand with the company's email infrastructure to provide a higher degree of fine-tuning, with the ability, for example, to choose between aggressive anti-spam measures (letting fewer emails in), and looser measures (letting more emails in).

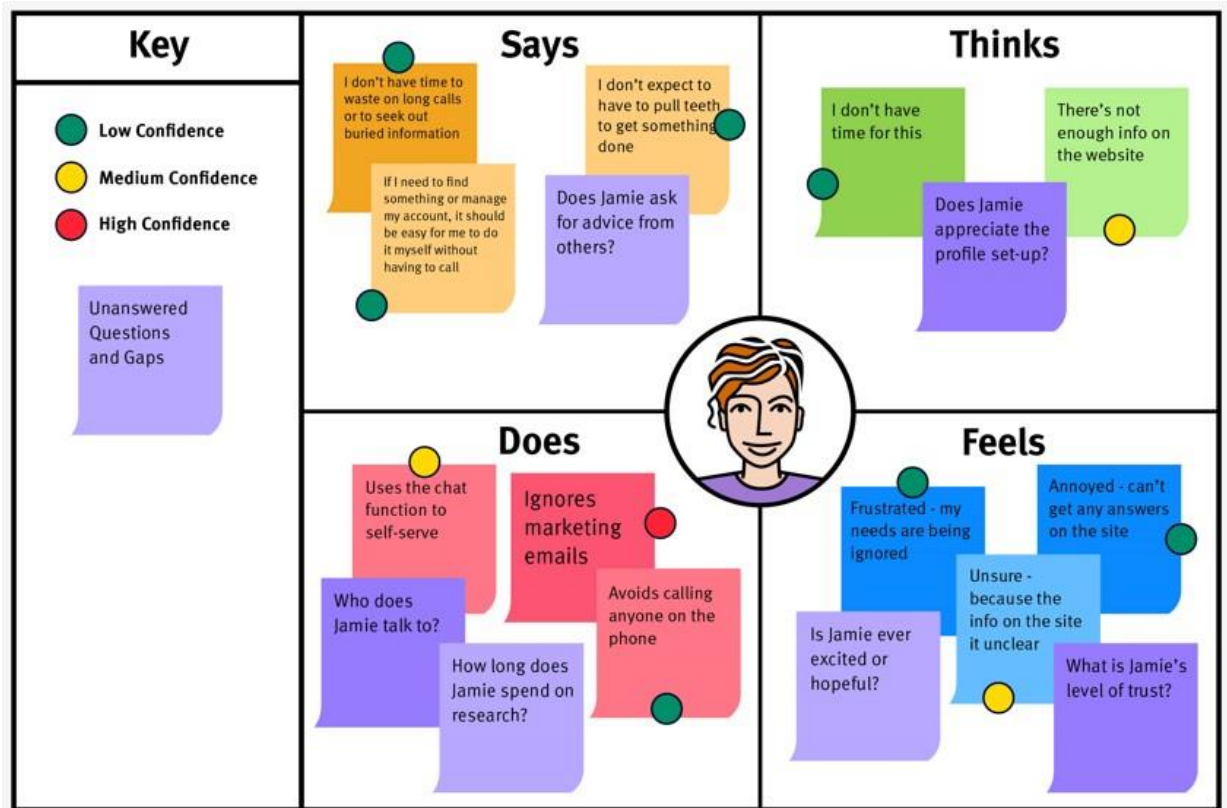
Companies that employ third party spam filters also benefit from more robust reporting and analytics so they can better understand email and spam activity within the company. This can be useful for proactively thwarting targeted phishing attacks against a company.

3. Desktop spam filters

Desktop spam filters are pieces of software that you purchase and install on your computer. They monitor messages once they make it past your email service provider's spam filters and add an extra layer of protection. The level of protection can be configured to specifically suit your spam tolerance and tastes. As the filters built into email service providers become more and more sophisticated, desktop spam filters have become less popular.

2. PROBLEM DEFINITION AND DESIGN THINKING

EMPATHY MAP



IDEATION AND BRAINSTROMING :



3. RESULT :



4. ADVANTAGES AND DISADVANTAGES

Study ID	Advantages	Disadvantages
S1 [1]	Combination of machine learning algorithms with user generated features	Users need to select features manually
S2 [2]		No classification algorithm is used
S3 [3]	Classification based on two algorithms	Threshold selection
S4 [5]	Combination of client and server side algorithms	Challenge-response technique suffers from server side traffic and user interaction problems
S5 [10]	Accurate as Naïve Bayesian with necessary feature extraction	Complex implementation
S6 [11]	Incorporating Apriori Algorithm	
S7 [12]	Used a weighting Mechanism to reduce false negatives	
S8 [7]		Suffers from implementation complexity
S9 [13]	Combination of two formulas gives better result in terms of false positives	
S10 [14]	Concludes SVM outperforms other algorithms and created a baseline for further comparison	
S11 [15]	Although SVM gives better results in Spam identification Bayesian is more feasible for mobile applications	Extensive feature engineering is needed for better accuracy
S12 [16]	Demonstrates the need of spam filtering in spite of having established email spam filtering	
S14 [18]	Lightweight and focuses on runtime	

5. APPLICATIONS :

A spam filtering solution cannot be 100 percent effective. However, a business email system without spam filtering is highly vulnerable, if not unusable. It is important to stop as much spam as you can, to protect your network from the many possible risks: viruses, phishing attacks, compromised web links and other malicious content.

Spam filters also protect your servers from being overloaded with non-essential emails, and the worse problem of being infected with spam software that may turn them into spam servers themselves.

By preventing spam email from reaching your employees' mailboxes, spam filters give an additional layer of protection to your users, your network, and your business.

6. CONCLUSION

In the last two decades, spam detection and filtration gained the attention of a sizeable research community. The reason for a lot of research in this area is its costly and massive effect in many situations like consumer behavior and fake reviews. The survey covers various machine learning techniques and models that the various researchers have proposed to detect and filter spam in emails and IoT platforms. The study categorized them as supervised, unsupervised, reinforcement learning, etc. The study compares these approaches and provides a summary of learned lessons from each category. This study concludes that most of the proposed email and IoT spam detection methods are based on supervised machine learning techniques. A labeled dataset for the supervised model training is a crucial and time-consuming task. Supervised learning algorithms SVM and Naïve Bayes outperform other models in spam detection. The study provides comprehensive insights of these algorithms and some future research directions for email spam detection and filtering.

7. FUTURE SCOPE

Efficient pattern detection in spam mail filtering plays crucial role. Using RFD model spam detection gives the spam patterns, non –spam patterns and general patterns which easily identify the whether the mail is spam or ham. The current method which uses the pattern detection method does not include the general patterns. RFD gives the general patterns of which user can decide to determine whether he wants to put the mail as spam or non-spam to avoid the loss of important mails. The images which are in forms of spams are also detected using File Properties, Histogram and Hough Transform. The current proposed system is for English language mails but as future scope we can design the system for multiple languages.

8. APPENDIX

```
9. import numpy as np
10.     import pandas as pd
11.     import matplotlib.pyplot as plt
12.     import tensorflow as tf
13.
14.     # For splitting dataset
15.     from sklearn.model_selection import train_test_split
16.
17.     # For preprocessing data
18.     import re
19.     import nltk
20.     import string
21.     from tensorflow.keras.preprocessing.text import
        Tokenizer
22.     from tensorflow.keras.preprocessing.sequence import
        pad_sequences
23.
24.     # For building the model
25.     from tensorflow.keras.callbacks import EarlyStopping
26.     from tensorflow.keras.models import Sequential
27.     from tensorflow.keras.layers import (
28.         Dense, Embedding, Dropout,
29.         GlobalAveragePooling1D, Input
30.     )
31.     from tensorflow.keras.models import Model
32.
33.     SEED = 42
34.     tf.random.set_seed(SEED)
35.     np.random.seed(SEED)
36.
37.     data_path = "spam.csv"
38.
39.     df = pd.read_csv(data_path, encoding='ISO-8859-1')
40.     df = df[['v1', 'v2']]
```

```

41.     df.columns = ['Category', 'Content']
42.     df.head()
43.
44.     plt.figure(figsize=(8, 6))
45.     df.Category.value_counts().plot(kind='bar',
color='green')
46.
47.     stopwords = nltk.corpus.stopwords.words('english')
48.     ps = nltk.PorterStemmer()
49.
50.     def clean_text(text):
51.         # Remove punctuation from text
52.         text = "".join([word.lower() for word in text if
word not in string.punctuation])
53.         tokens = re.split('\W+', text)
54.         # Get Origin of words (books -> book)
55.         text = " ".join([ps.stem(word) for word in tokens
if word not in stopwords])
56.         return text
57.
58.         df['Clean_text'] = df['Content'].apply(lambda x:
clean_text(x))
59.     df.head()
60.
61.     # {ham:0, spam:1}
62.     df['label'] =
df['Category'].astype('category').cat.codes
63.     df.head()
64.
65.     X_train,X_test,y_train,y_test=train_test_split(
66.         df['Clean_text'],df['label'],test_size=0.25,random_state=SEE
D)
67.
68.     max_len=50 # max words in every input string

```

```

69.     trunc_type='post' # if number of words are more than
        max_len then trunc from the last of the string
70.     padding_type='post' # if number of words are less than
        max_len then add
71.     oov_token_1='<OOV>' # out of vocabulary token
72.     vocab_size=500
73.
74.     tokenizer=Tokenizer(num_words=vocab_size,char_level=False,oov_token=oov_token_1)
75.     tokenizer.fit_on_texts(X_train)
76.
77.     print(f"Total token : {len(tokenizer.word_index)}")
78.     for token_word in tokenizer.word_index:
79.         print(f"{token_word}:
            {tokenizer.word_index[token_word]}")
80.         if tokenizer.word_index[token_word] >10:
81.             break
82.
83.     X_train = tokenizer.texts_to_sequences(X_train)
84.     X_train =
        pad_sequences(X_train,maxlen=max_len,padding=padding_type,truncating=trunc_type)
85.     X_test=tokenizer.texts_to_sequences(X_test)
86.     X_test=pad_sequences(X_test,maxlen=max_len,padding=padding_type,truncating=trunc_type)
87.     X_train.shape, X_test.shape
88.
89.     input_shape = X_train.shape[1:]
90.     input_shape
91.
92.     input = Input(shape=input_shape)
93.     x = Embedding(input_dim=vocab_size,
        output_dim=12)(input)
94.     x = GlobalAveragePooling1D()(x)
95.     x = Dense(24, activation='relu')(x)
96.     x = Dropout(0.2)(x)

```

```
97.     output = Dense(1, activation='sigmoid')(x)
98.     model = Model(inputs=input, outputs=output)
99.     model.compile(optimizer = 'adam', loss =
    'binary_crossentropy', metrics=['accuracy'])
100.
101.     model.summary()
102.
103.     tf.keras.utils.plot_model(model)
104.
105.     num_epochs=50
106.     # If val_loss increases five times it will stop model
    fitting
107.     early_stop=EarlyStopping(monitor='val_loss',patience=5)
108.     history=model.fit(
109.         X_train,
110.         y_train,
111.         epochs=num_epochs,
112.         validation_data=(X_test,y_test),
113.         callbacks=[early_stop]
114.     )
115.
116.     model.evaluate(X_test, y_test)
117.
118.     def plot_loss_curves(history):
119.         loss = history.history['loss']
120.         val_loss = history.history['val_loss']
121.
122.         accuracy = history.history['accuracy']
123.         val_accuracy = history.history['val_accuracy']
124.
125.         epochs = range(len(history.history['loss']))
126.
127.         # Plot loss
128.         plt.plot(epochs, loss, label='training_loss')
129.         plt.plot(epochs, val_loss, label='val_loss')
130.         plt.title('Loss')
```

```
131.     plt.xlabel('Epochs')
132.     plt.legend()
133.
134.     # Plot accuracy
135.     plt.figure()
136.     plt.plot(epochs, accuracy, label='training_accuracy')
137.     plt.plot(epochs, val_accuracy, label='val_accuracy')
138.     plt.title('Accuracy')
139.     plt.xlabel('Epochs')
140.     plt.legend();
141.
142.     plot_loss_curves(history)
143.
```